

# No Regret Bound for Extreme Bandits

Robert Nishihara  
UC Berkeley  
rkn@eecs.berkeley.edu

David Lopez-Paz  
Facebook AI Research  
dlp@fb.com

Léon Bottou  
Facebook AI Research  
leonb@fb.com

## Abstract

Algorithms for hyperparameter optimization abound, all of which work well under different and often unverifiable assumptions. Motivated by the general challenge of sequentially choosing which algorithm to use, we study the more specific task of choosing among distributions to use for random hyperparameter optimization. This work is naturally framed in the extreme bandit setting, which deals with sequentially choosing which distribution from a collection to sample in order to minimize (maximize) the single best cost (reward). Whereas the distributions in the standard bandit setting are primarily characterized by their means, a number of subtleties arise when we care about the minimal cost as opposed to the average cost. For example, there may not be a well-defined “best” distribution as there is in the standard bandit setting. The best distribution depends on the rewards that have been obtained and on the remaining time horizon. Whereas in the standard bandit setting, it is sensible to compare policies with an oracle which plays the single best arm, in the extreme bandit setting, there are multiple sensible oracle models. We define a sensible notion of “extreme regret” in the extreme bandit setting, which parallels the concept of regret in the standard bandit setting. We then prove that no policy can asymptotically achieve no extreme regret.

## 1 Introduction

Our motivation comes from hyperparameter optimization and more generally from the challenge of mini-

mizing a black-box objective  $f: \Omega \rightarrow [0, 1]$  which we can only evaluate pointwise. As an example,  $\omega \in \Omega$  may parameterize the architecture of a convolutional network, and  $f(\omega)$  may be the validation error when the network with that architecture is trained on a particular data set. A number of approaches have been applied to the optimization of  $f$  including Bayesian optimization, covariance matrix adaptation, random search, and a variety of other methods (for an incomplete list, see Bergstra and Bengio (2012); Bergstra et al. (2011); Snoek et al. (2012); Hansen (2006); Wang et al. (2013); Lagarias et al. (1998); Powell (2006); Duchi et al. (2015)).

In some sense, random search is the benchmark of choice. Whereas other approaches work well under various and often unverifiable conditions (such as smoothness or convexity of the objective), random search has strong finite-sample guarantees that hold without any assumptions on the function under consideration. This guarantee is illustrated by the so-called *rule of 59*,<sup>1</sup> which states that the best of 59 random samples will be in the best 5 percent of all samples with probability at least 0.95. More generally, any distribution over the set of hyperparameters  $\Omega$  induces a distribution  $\mu$  over the validation error in  $[0, 1]$ . Let  $F_\mu$  be the cumulative distribution function of  $\mu$ , and suppose that  $F_\mu$  is continuous. Suppose that  $x_1, \dots, x_T$  are independent and identically-distributed samples from  $\mu$  (obtained, for instance, by independently sampling hyperparameters  $\omega_t$  and evaluating  $x_t = f(\omega_t)$  for  $1 \leq t \leq T$ ). The following is known.

**Lemma 1.** *The distribution of the extreme cost  $\min\{x_1, \dots, x_T\}$  is easily described with quantiles. We have  $P(F_\mu(\min\{x_1, \dots, x_T\}) \leq \alpha) = 1 - (1 - \alpha)^T$ . More specifically,  $F_\mu(\min\{x_1, \dots, x_T\})$  is a Beta(1, T) random variable.*

*Proof.* The event  $F_\mu(\min\{x_1, \dots, x_T\}) > \alpha$  happens if and only if  $F_\mu(x_t) > \alpha$  for each  $t$ , which happens

Appearing in Proceedings of the 19<sup>th</sup> International Conference on Artificial Intelligence and Statistics (AISTATS) 2016, Cadiz, Spain. JMLR: W&CP volume 41. Copyright 2016 by the authors.

<sup>1</sup>Though they are known, the rule of 59 and Lemma 1 do not appear in Bergstra and Bengio (2012), and they are difficult to find in the literature.

with probability  $(1 - \alpha)^T$ . Differentiating the resulting cumulative distribution function gives the density function of a Beta(1,  $T$ ) random variable.  $\square$

The generality of Lemma 1 comes at a price. The guarantee is given with respect to the distribution  $\mu$ , but there is no guarantee about  $\mu$  itself. Different induced distributions  $\mu$  may arise from different parameterizations of the hyperparameter space  $\Omega$  (for example, from the decision to put a uniform or a log-uniform distribution over a coordinate of  $\omega$ ), and the allocation of mass over  $[0, 1]$  may vary wildly based on these choices.

Furthermore, the flip side of making no assumptions on the underlying objective is that random search fails to adapt to easy problems. When the objective under consideration satisfies various regularity conditions (as real-world objectives often do), more heavily-engineered approaches will likely outperform random search. That said, it is not clear how to know that a given algorithm is outperforming random search without also running random search. For this reason, the benefits of a potentially faster algorithm are blunted when one must also run the slow algorithm to verify the performance of the fast algorithm.

Given the variety of existing hyperparameter optimization algorithms, it would be desirable to devise a strategy for sequentially choosing which algorithm to use in a way that performs nearly as well as if we had only used the single best algorithm. We consider the simpler problem of choosing which of several distributions over hyperparameters to use for random search. In Theorem 11, we show that even in this simplified setting, no strategy guarantees performance that is asymptotically as good as the single best distribution, at least not without stronger assumptions.

We will frame our negative result in the extreme bandit setting (Carpentier and Valko, 2014), also called the max  $K$ -armed bandit setting (Cicirello and Smith, 2005). Prior work has focused on designing algorithms that perform asymptotically as well as the single best distribution under parametric (or semiparametric) assumptions on the possible distributions (Cicirello and Smith, 2005; Carpentier and Valko, 2014). Instead, we focus on probing the difficulty of the problem, pointing out a number of subtleties that arise in this setting that do not show up in the conventional bandit setting.

## 2 The Extreme Bandit Setting

Cicirello and Smith (2005) introduce the extreme bandit problem (they call it the max  $K$ -armed bandit problem) as follows. We are given a tuple of unknown

distributions (arms)  $\mu_1^K = (\mu_1, \dots, \mu_K)$ . The  $k$ th distribution generates sample  $x_{k,t}$  at time  $t$ , for integer  $t \geq 1$ , and all of the samples  $x_{k,t}$  are independent. A policy  $\pi$  is a function that, at each time  $t$ , chooses the index  $k_t$  of a distribution to sample based on the previously observed samples. That is,

$$k_t = \pi(\underbrace{k_1, \dots, k_{t-1}}_{\text{past arm choices}}, \underbrace{x_{k_1,1}, \dots, x_{k_{t-1},t-1}}_{\text{past values}}).$$

We would like to compare the performance of a policy  $\pi$  to that of an oracle policy  $\pi_*$  that has access to knowledge of the distributions  $\mu_1^K$ , so

$$k_t^* = \pi_*(\mu_1^K, k_1^*, \dots, k_{t-1}^*, x_{k_1^*,1}, \dots, x_{k_{t-1}^*,t-1}).$$

Both Cicirello and Smith (2005) and Carpentier and Valko (2014) phrase their results in terms of the maximization of a reward rather than the minimization of a cost. They define the “regret” of policy  $\pi$  with respect to the oracle  $\pi_*$  over a time horizon of  $T$  as

$$G_T^{\pi, \pi_*} = \mathbb{E} \left[ \max_{t \leq T} x_{k_t^*, t} \right] - \mathbb{E} \left[ \max_{t \leq T} x_{k_t, t} \right].$$

Under semiparametric assumptions on  $\mu_1^K$ , Carpentier and Valko (2014) exhibit a policy  $\pi$  such that

$$G_T^{\pi, \pi_*} \text{ is } o \left( \mathbb{E} \left[ \max_{t \leq T} x_{k_t^*, t} \right] \right) \quad (1)$$

or equivalently,

$$\lim_{T \rightarrow \infty} \frac{\mathbb{E} [\max_{t \leq T} x_{k_t, t}]}{\mathbb{E} [\max_{t \leq T} x_{k_t^*, t}]} \rightarrow 1. \quad (2)$$

The result in Equation 1 is superficially similar to results in the standard bandit setting. However, while the condition in Equation 1 is sensible for the setting considered by Carpentier and Valko (2014) (where the distributions  $\mu_1^K$  have unbounded support), it is particularly sensitive to the nature of the distributions. For instance, the result in Equation 1 is trivially achieved when the distributions have bounded support (for example, when the support is contained in  $[0, 1]$  as in hyperparameter optimization). In this case, both the numerator and denominator converge to the upper bound of the support and  $G_T^{\pi, \pi_*} \rightarrow 0$  (for any policy that chooses each distribution infinitely often).

Furthermore, the condition in Equation 2 is asymmetric with respect to maximization and minimization. When performing minimization of a cost instead of maximization of a reward (using distributions supported in  $[0, 1]$ ), both  $\mathbb{E} [\min_{t \leq T} x_{k_t, t}]$  and  $\mathbb{E} [\min_{t \leq T} x_{k_t^*, t}]$  may approach 0, in which case the

ratio may exhibit radically different behavior. In Example 2 and Example 3, we demonstrate some of the peculiarities of this performance metric in the minimization setting.

**Example 2.** Suppose  $\mu_1$  is a Bernoulli distribution with mean parameter  $0 < p < 1$  and suppose that  $\mu_2$  is a point mass on 1. Consider a policy  $\pi$  which chooses  $\mu_2$  at  $t = 1$  and then chooses  $\mu_1$  for all  $t > 1$  and a policy  $\pi_*$  which always chooses  $\mu_1$ . We have

$$\lim_{T \rightarrow \infty} \frac{\mathbb{E}[\min_{t \leq T} x_{k_t, t}]}{\mathbb{E}[\min_{t \leq T} x_{k_t^*, t}]} = \lim_{T \rightarrow \infty} \frac{p^{T-1}}{p^T} = \frac{1}{p},$$

which remains bounded away from 1 even though the policy  $\pi$  acted optimally at every time step after  $t = 1$ .

**Example 3.** Suppose  $\mu_1$  is the uniform distribution over  $[0, 1]$  and suppose that  $\mu_2$  is a point mass on 1. Consider a policy  $\pi$  which chooses  $\mu_2$  at  $t = 1$  and then chooses  $\mu_1$  for all  $t > 1$  and a policy  $\pi_*$  which always chooses  $\mu_1$ . We have

$$\lim_{T \rightarrow \infty} \frac{\mathbb{E}[\min_{t \leq T} x_{k_t, t}]}{\mathbb{E}[\min_{t \leq T} x_{k_t^*, t}]} = \lim_{T \rightarrow \infty} \frac{T^{-1}}{(T+1)^{-1}} \rightarrow 1.$$

Note above that the minimum of  $T$  independent uniform random variables is a Beta(1,  $T$ ) random variable, which has mean  $1/(T+1)$ .

Despite the fact that the policy  $\pi$  acts optimally at every time step other than  $t = 1$  in both Example 2 and Example 3, the ratios of their expectations to that of the oracle  $\pi_*$  exhibit wildly different behaviors.

To avoid this sensitivity, we define “extreme regret” as follows.

**Definition 4.** We define the extreme regret of the policy  $\pi$  with respect to the oracle policy  $\pi_*$  over a time horizon of  $T$  as

$$R_T^{\pi, \pi_*} = \frac{1}{T} \min_{T' \geq 1} \left\{ T' : \mathbb{E} \left[ \min_{t \leq T'} x_{k_t, t} \right] \leq \mathbb{E} \left[ \min_{t \leq T} x_{k_t^*, t} \right] \right\}.$$

Note that  $R_T^{\pi, \pi_*}$  depends on the tuple of distributions  $\mu_1^K$ , but we suppress this dependence in our notation.

Then  $R_T^{\pi, \pi_*}$  is essentially the ratio of the time horizons  $T'$  to  $T$  over which the policy  $\pi$  and the oracle  $\pi_*$  perform equally well. This definition is sensible regardless of whether the samples are bounded or unbounded, whether we care about minimization or maximization, and regardless of how we scale or translate the distributions. Note that in both Example 2 and Example 3, we have  $R_T^{\pi, \pi_*} = \frac{T+1}{T} \rightarrow 1$ . Despite its apparent difference, as we discuss in Section 2.1, Definition 4 is closely related to the notion of regret used in the standard bandit setting.

**Definition 5.** We say that policy  $\pi$  achieves “no extreme regret” with respect to the oracle  $\pi_*$  if  $\limsup_T R_T^{\pi, \pi_*} \leq 1$  for all tuples of distributions  $\mu_1^K$ .

Definition 5 is fairly lenient. Had we defined “no extreme regret” using the condition given in Equation 1, our main result in Theorem 11 could have been made even stronger, but we view that as undesirable as illustrated by Example 2 and Example 3. Moreover, the quantities in Definition 4 and Definition 5 closely parallel quantities of interest in the standard bandit setting, as we show in Section 2.1.

## 2.1 Analogy with the Standard Bandit Setting

Definition 4 and Definition 5 parallel the intuition of the standard bandit setting, which (when minimizing a cost) studies the rate of convergence of

$$\frac{\mathbb{E} \left[ \sum_{t=1}^T x_{k_t, t} \right] - \min_k \mathbb{E} \left[ \sum_{t=1}^T x_{k, t} \right]}{\min_k \mathbb{E} \left[ \sum_{t=1}^T x_{k, t} \right]} \rightarrow 0. \quad (3)$$

Adding 1 to both sides, this is the same as studying the rate of convergence of

$$\frac{\mathbb{E} \left[ \sum_{t=1}^T x_{k_t, t} \right]}{T \min_k \mathbb{E}[x_{k, t}]} \rightarrow 1.$$

Now, observe that we have

$$\begin{aligned} & \frac{\mathbb{E} \left[ \sum_{t=1}^T x_{k_t, t} \right]}{T \min_k \mathbb{E}[x_{k, t}]} \\ & \approx \frac{1}{T} \min_{T' \geq 1} \left\{ T' : \frac{\mathbb{E} \left[ \sum_{t=1}^{T'} x_{k_t, t} \right]}{\min_k \mathbb{E}[x_{k, t}]} \leq T' \right\} \\ & = \frac{1}{T} \min_{T' \geq 1} \left\{ T' : \mathbb{E} \left[ \sum_{t=1}^{T'} x_{k_t, t} \right] \leq \min_k \mathbb{E} \left[ \sum_{t=1}^{T'} x_{k, t} \right] \right\}, \end{aligned} \quad (4)$$

which is essentially the ratio of the time horizons over which the policy and the oracle perform equally well. The two sides of the approximate equality in Equation 4 differ by at most  $1/T$ . In the standard bandit setting, the term “regret” often refers to the numerator in Equation 3 and not the quantity in Equation 4. However, as the above computation shows, the two quantities are closely related, and they capture the same phenomenon. We will phrase our results in terms of the quantity  $R_T^{\pi, \pi_*}$  from Definition 4, which parallels the quantity in Equation 4.

## 3 Oracle Models

In the standard multi-armed bandit setting, if an oracle with knowledge of the distributions of the arms

seeks to minimize the expected sum of the losses, it should simply choose to play the arm with the lowest mean. This is true regardless of the time horizon. By analogy with the usual multi-armed bandit setting, Cicirello and Smith (2005) and Carpentier and Valko (2014) both consider the oracle policy in Definition 6 that plays the single “best” arm.

**Definition 6** (single-armed oracle). *The single-armed oracle is the oracle, which over a time horizon of  $T$ , plays the single best arm*

$$\arg \min_k \mathbb{E} \left[ \min_{t \leq T} x_{k,t} \right].$$

The single-armed oracle provides a good benchmark for comparison, but it is not the optimal oracle policy. When the time horizon is known in advance, the optimal oracle policy is given in Definition 7.

**Definition 7** (optimal oracle). *The optimal oracle over a time horizon of  $T$  plays the policy that solves*

$$\arg \min_{\pi} \mathbb{E} \left[ \min_{t \leq T} x_{k_t, t} \right].$$

When the time horizon is not known in advance, one possible oracle strategy is a greedy strategy given in Definition 8.

**Definition 8** (greedy oracle). *The greedy oracle chooses the arm  $k_t^*$  at time  $t$  that gives the maximal expected improvement over the current best value  $y_{t-1} = \min_{s \leq t-1} x_{k_s^*, s}$ . That is,*

$$k_t^* = \arg \min_k \mathbb{E} \left[ \min \{x_{k,t}, y_{t-1}\} \mid x_{k_1^*, 1}, \dots, x_{k_{t-1}^*, t-1} \right].$$

Unlike the greedy oracle, both the single-armed oracle and the optimal oracle require knowledge of the time horizon. Indeed, as shown in Example 9, the notion of a “best” arm is not well-defined outside of a specific time horizon. The best arm depends on the time horizon. This point contrasts sharply with the usual multi-armed bandit setting.

**Example 9.** *Suppose we have an infinite collection of arms  $\mu_s$  indexed by  $0 < s < 1$ . Let  $x_{s,t}$  be a sample from  $\mu_s$  and suppose that  $P(x_{s,t} = s) = s$  and  $P(x_{s,t} = 1) = 1 - s$ . Then the optimal  $s$  is  $\Theta((\log T)/T)$ .*

We elaborate on Example 9 in Appendix A. One difference between the single-armed oracle and the optimal oracle is that the optimal oracle can adapt its strategy based on the samples that it receives, whereas the single-armed oracle is non-adaptive. Its strategy is fixed ahead of time. Example 10 shows that the single-armed oracle is not even the optimal non-adaptive oracle. A mixed strategy may outperform any policy that plays only a single arm.

**Example 10.** *Consider a time horizon  $T = 2$  and consider two arms. Suppose that samples  $x_{1,t}$  from  $\mu_1$  deterministically equal  $1/2$  and that samples  $x_{2,t}$  from  $\mu_2$  satisfy  $P(x_{2,t} = 0) = 1/4$  and  $P(x_{2,t} = 1) = 3/4$ . Then*

$$\begin{aligned} \mathbb{E} \min_{1 \leq t \leq 2} x_{1,t} &= \frac{1}{2} \\ \mathbb{E} \min_{1 \leq t \leq 2} x_{2,t} &= \frac{9}{16} \\ \mathbb{E} \min \{x_{1,1}, x_{2,2}\} &= \frac{3}{8}. \end{aligned}$$

*This example shows that a fixed strategy that plays both arms can outperform any policy that plays a single-arm.*

We described three different oracle models above. One caveat is that in the event that there is a well-defined best arm, that is, some arm  $k_*$  such that  $P(x_{k_*, t} \leq \alpha) \geq P(x_{k, t} \leq \alpha)$  for all  $k$  and all  $0 \leq \alpha \leq 1$ , then these three oracles all coincide and we need not worry about which oracle to use for comparison. This is roughly the case in prior work. Cicirello and Smith (2005) and Carpentier and Valko (2014) make (semi)parametric assumptions on the distributions of the arms which essentially restrict the setting to have a well-defined best arm.

Despite the fact that the single-armed oracle is not the optimal oracle strategy, it is often a sufficiently strong baseline for measuring the performance of our policies. When we cannot even do as well as the single-armed oracle, as will be the case in Theorem 11, then we also cannot do as well as the optimal oracle. For the remainder of the paper, we will compare to the single-armed oracle. However, the results necessarily hold for comparisons to the optimal oracle as well.

## 4 Main Result

Theorem 11 shows that no policy can be guaranteed to perform asymptotically as well as the single best distribution. That is, it is impossible to achieve “no extreme regret” in the extreme bandit problem. This result contrasts sharply with results in the standard bandit setting, where it is possible to achieve no regret under relatively mild conditions on the distributions  $\mu_1^K = (\mu_1, \dots, \mu_K)$ .

**Theorem 11.** *For any policy  $\pi$ , there exist distributions  $\mu_1^K$  such that  $\limsup_T R_T^{\pi, \pi_*} \geq K$ , where  $\pi_*$  is the single-armed oracle.*

We prove Theorem 11 in Section 4.3. The main components of the proof are Lemma 13, which upper bounds the performance of the single-armed oracle and

Lemma 15, which lower bounds the performance of the policy  $\pi$ .

This result shows that the extreme bandit problem is fundamentally different from the standard multi-armed bandit problem, where a variety of policies perform asymptotically as well as the single best arm. Indeed, in the standard bandit problem, the arms are primarily characterized by their means, and so it suffices to estimate the means of the arms and play the best one. However, as discussed in Example 9, there is no well-defined best arm in the extreme bandit problem. Our construction will create a situation where the “best” arm periodically switches among the  $K$  distributions so that the policy  $\pi$  often ends up choosing the “wrong” arm.

For  $i \geq 1$ , let  $\alpha_i = (8K)^{-(i!)^2}$ . Our construction will involve a sum of point masses at the values  $\alpha_i$ . It is easily verified that the sequence  $\alpha_i$  satisfies the conditions in Lemma 12.

**Lemma 12.** *The sequence  $\alpha_i$  satisfies the following properties.*

$$(A) \sum_{j=1}^{\infty} \alpha_j \leq 1/2$$

$$(B) \alpha_i \leq \frac{1}{4(1+i)}$$

$$(C) \sum_{j=i+1}^{\infty} \alpha_j \leq \frac{\alpha_i}{iK}$$

$$(D) \alpha_i \leq \alpha_{i-1}^i 2^{-i}.$$

Henceforth, we will not need the exact values of the sequence, we will only need the properties enumerated in Lemma 12. For  $b = (b_1, b_2, \dots) \in \{1, \dots, K\}^\infty$ , define the tuple of distributions  $\mu_1^K(b) = (\mu_1(b), \dots, \mu_K(b))$  via

$$\mu_k(b) = \gamma_k(b) \delta_1 + \sum_{i=1}^{\infty} \mathbb{1}[b_i = k] \alpha_i \delta_{\alpha_i}$$

where

$$\gamma_k(b) = 1 - \sum_{i=1}^{\infty} \mathbb{1}[b_i = k] \alpha_i.$$

Here,  $\delta_c$  represents a point mass at  $c$ ,  $\mathbb{1}[\xi]$  is the  $\{0, 1\}$ -indicator function of the event  $\xi$ , and  $\gamma_k(b)$  is chosen to make  $\mu_k(b)$  a probability measure. Let  $M_K$  be the set of tuples of distributions that can be obtained in this way. The value  $b_i$  simply assigns the point mass  $\delta_{\alpha_i}$  to one of the  $K$  distributions. We let  $D$  denote the distribution over the set  $\{1, \dots, K\}^\infty$  defined so that the  $b_i$ ’s are independent uniform random variables in  $\{1, \dots, K\}$ .

Define the time horizon  $T_i = \lceil \log(1/\alpha_i)/\alpha_i \rceil$ . Instead of controlling  $R_T^{\pi, \pi^*}$  for every  $T$ , we will control the quantity specifically for the time horizons  $T_i$ . In our construction, the  $b_i$ th arm in the tuple will be the best

arm over the time horizon  $T_i$ , and the other arms will be substantially worse. We will show that, for a fixed  $i$ , we can construct a tuple  $\mu_1^K$  so that the policy  $\pi$  takes roughly  $K$  times longer than the single-armed oracle  $\pi_*$  to obtain the value  $\alpha_i$  (that is,  $\pi_*$  requires roughly  $T_i$  samples and  $\pi$  requires roughly  $T_i' \approx KT_i$  samples). Using the probabilistic method, we will then show that we can find a tuple  $\mu_1^K$  so that the policy takes roughly  $K$  times longer than the oracle to obtain the value  $\alpha_i$  for infinitely many values of  $i$ .

#### 4.1 Upper Bound on Oracle Performance

We begin by giving an upper bound on the performance of the oracle policy that plays the single best arm over the time horizon  $T_i$ . This bound holds uniformly over  $M_K$ .

**Lemma 13.** *Suppose that  $\mu_1^K(b) \in M_K$ . If  $\pi_*$  is the single-armed oracle from Definition 6, then*

$$\mathbb{E} \left[ \min_{t \leq T_i} x_{k^*, t} \right] < 2\alpha_i.$$

*Proof.* Recall that  $b_i$  is the index of the distribution that has a point mass at  $\alpha_i$ . We have

$$\mathbb{E} \left[ \min_{t \leq T_i} x_{k^*, t} \right] = \min_k \mathbb{E} \left[ \min_{t \leq T_i} x_{k, t} \right] \leq \mathbb{E} \left[ \min_{t \leq T_i} x_{b_i, t} \right].$$

The term on the right hand side can be rewritten as

$$\begin{aligned} & \mathbb{E} \left[ \mathbb{1} \left[ \min_{t \leq T_i} x_{b_i, t} \leq \alpha_i \right] \min_{t \leq T_i} x_{b_i, t} \right] \\ & + \mathbb{E} \left[ \mathbb{1} \left[ \min_{t \leq T_i} x_{b_i, t} > \alpha_i \right] \min_{t \leq T_i} x_{b_i, t} \right] \\ & \leq \alpha_i P \left[ \min_{t \leq T_i} x_{b_i, t} \leq \alpha_i \right] + P \left[ \min_{t \leq T_i} x_{b_i, t} > \alpha_i \right] \\ & \leq \alpha_i + P \left[ \min_{t \leq T_i} x_{b_i, t} > \alpha_i \right]. \end{aligned}$$

The first inequality follows by upperbounding the term  $\min_{t \leq T_i} x_{b_i, t}$  by  $\alpha_i$  in the first term and by 1 in the second term. The second inequality follows by upperbounding the first probability by 1. To finish the lemma, note that

$$P \left[ \min_{t \leq T_i} x_{b_i, t} > \alpha_i \right] \leq (1 - \alpha_i)^{T_i} < e^{-\alpha_i T_i} \leq \alpha_i,$$

where the third inequality uses the definition  $T_i = \lceil \log(1/\alpha_i)/\alpha_i \rceil$ .  $\square$

#### 4.2 Lower Bound on Performance of $\pi$

Here, we give a lower bound on the performance of any fixed policy  $\pi$ , when averaged over a collection of tuples of distributions.



Define the time horizon  $T'_i = \lfloor c_i K \log(1/\alpha_i)/\alpha_i \rfloor$ , where  $c_i = (1 - 1/i)/((1 + 1/i)^2 + 2/i)$ . The constant  $c_i$  is a correction term that converges to 1 as  $i \rightarrow \infty$ . Its specific value is not meaningful. The goal of this section is roughly to show that the performance of the policy  $\pi$  over a time horizon of  $T'_i$  is comparable to the performance of the oracle policy over a time horizon of  $T_i$ .

Throughout this section, we will fix an index  $i$  and we fix  $b_j$  for all  $j \neq i$ . Then we define the sequence  $b^{k'} = (b_1^{k'}, b_2^{k'}, \dots)$  via  $b_j^{k'} = b_j$  for  $j \neq i$  and  $b_i^{k'} = k'$ . The  $K$  tuples  $\mu_1^K(b^{k'})$  for different values of  $k'$  are identical in all respects except for the index of the distribution that possesses the point mass  $\delta_{\alpha_i}$  and the amount of mass  $\gamma_k(b^{k'})$  that the  $k$ th distribution in the  $k'$ th tuple assigns to  $\delta_1$ .

Define the tuple of distributions  $\eta_1^K(\bar{b}) = (\eta_1(\bar{b}), \dots, \eta_K(\bar{b}))$  by  $\eta_k(\bar{b}) = \frac{1}{K} \sum_{k'=1}^K \mu_k(b^{k'})$ . Let  $\gamma_k(\bar{b}) := \frac{1}{K} \sum_{k'=1}^K \gamma_k(b^{k'})$  denote the probability that  $\eta_k(\bar{b})$  assigns to the value 1. The tuple  $\eta_1^K(\bar{b})$  is the average of the tuples  $\mu_1^K(b^{k'})$  over the different values of  $k'$ .

We begin with Lemma 14 which compares the probability that policy  $\pi$  obtains the value  $\alpha_i$  when averaged over the tuples  $\mu_1^K(b^{k'})$  with the probability that  $\pi$  obtains the value  $\alpha_i$  in the tuple  $\eta_1^K(\bar{b})$ . This comparison is helpful because each distribution in the tuple  $\eta_1^K(\bar{b})$  assigns the same mass of  $\alpha_i/K$  to  $\alpha_i$  and so the probability that  $\pi$  obtains  $\alpha_i$  when run on the tuple  $\eta_1^K(\bar{b})$  does not depend on  $\pi$  (it is simply  $(1 - \alpha_i/K)^{T'_i}$  where  $T'_i$  is the time horizon). Of course, as stated, we are actually concerned with the probability that  $\pi$  obtains a value less than or equal to  $\alpha_i$ , but because of Lemma 12(C), the contribution of the smaller terms will not be too great.

**Lemma 14.** *We have*

$$\begin{aligned} & \frac{1}{K} \sum_{k'=1}^K P \left[ \min_{t \leq T'_i} x_{k_t, t} \geq \alpha_{i-1} \mid \mu_1^K(b^{k'}) \right] \\ & \geq cP \left[ \min_{t \leq T'_i} x_{k_t, t} \geq \alpha_{i-1} \mid \eta_1^K(\bar{b}) \right], \end{aligned}$$

where  $c = e^{-\frac{2\alpha_i T'_i}{iK}}$ . In our notation, we condition on  $\mu_1^K(b^{k'})$  to indicate the tuple of distributions being used.

*Proof.* Define  $S(\pi, \mu_1^K, T)$  to be the set of actions and values that can be obtained by following policy  $\pi$  on the tuple  $\mu_1^K$  for a time horizon of  $T$ . That is,

$$\begin{aligned} & S(\pi, \mu_1^K, T) \\ & = \left\{ (k_t, x_t)_{t=1}^T : \begin{array}{l} k_t = \pi(k_1, \dots, k_{t-1}, x_1, \dots, x_{t-1}) \\ x_t \in \text{supp}(\mu_{k_t}) \end{array} \right\}, \end{aligned}$$

where  $\text{supp}(\mu_{k_t})$  is the support of the distribution  $\mu_{k_t}$ . Then define  $S(\pi, \mu_1^K, T, i)$  to be the subset of  $S(\pi, \mu_1^K, T)$  such that all values are greater than or equal to  $\alpha_{i-1}$ . That is,

$$S(\pi, \mu_1^K, T, i) = \{(k_t, x_t)_{t=1}^T \in S(\pi, \mu_1^K, T) : x_t \geq \alpha_{i-1}\}.$$

Critically, note that

$$\begin{aligned} S(\pi, \eta_1^K(\bar{b}), T'_i, i) &= S(\pi, \mu_1^K(b^1), T'_i, i) \\ &\vdots \\ &= S(\pi, \mu_1^K(b^K), T'_i, i). \end{aligned} \tag{5}$$

Equation 5 holds because the supports of the tuples  $\mu_1^K(b^{k'})$  and  $\eta_1^K(\bar{b})$  only differ on  $\alpha_i$ , but we are considering only values that are at least  $\alpha_{i-1}$ , so this difference does not affect the sets. We shall refer to this common set as  $S$ . We have

$$\begin{aligned} & P \left[ \min_{t \leq T'_i} x_{k_t, t} \geq \alpha_{i-1} \mid \mu_1^K(b^{k'}) \right] \\ &= \sum_S \left( \prod_{j=1}^{i-1} \alpha_j^{|\{t: x_t = \alpha_j\}|} \prod_{k=1}^K \gamma_k(b^{k'})^{|\{t: k_t = k, x_t = 1\}|} \right). \end{aligned} \tag{6}$$

It follows that

$$\begin{aligned} & \frac{1}{K} \sum_{k'=1}^K P \left[ \min_{t \leq T'_i} x_{k_t, t} \geq \alpha_{i-1} \mid \mu_1^K(b^{k'}) \right] \\ &= \frac{1}{K} \sum_{k'=1}^K \sum_S \left( \prod_{j=1}^{i-1} \alpha_j^{|\{t: x_t = \alpha_j\}|} \prod_{k=1}^K \gamma_k(b^{k'})^{|\{t: k_t = k, x_t = 1\}|} \right) \\ &= \sum_S \left( \prod_{j=1}^{i-1} \alpha_j^{|\{t: x_t = \alpha_j\}|} \left( \frac{1}{K} \sum_{k'=1}^K \prod_{k=1}^K \gamma_k(b^{k'})^{|\{t: k_t = k, x_t = 1\}|} \right) \right), \end{aligned} \tag{7}$$

where the first equality uses Equation 6 and the second equality simply rearranges the terms. We would like to essentially apply Jensen's inequality to say something like

$$\frac{1}{K} \sum_{k'=1}^K \prod_{k=1}^K \gamma_k(b^{k'})^{|\{t: k_t = k, x_t = 1\}|} \geq \prod_{k=1}^K \gamma_k(\bar{b})^{|\{t: k_t = k, x_t = 1\}|}. \tag{9}$$

Unfortunately, despite the fact that  $\gamma_k$  is convex on the relevant region,  $\prod_{k=1}^K \gamma_k$  is not quite convex. However, it is nearly convex, and as we show in Lemma 17, Equation 9 holds up to a correction factor of  $e^{-\frac{2\alpha_i T'_i}{iK}}$ .

Using this result in Equation 8 gives

$$\begin{aligned}
 & \frac{1}{K} \sum_{k'=1}^K P \left[ \min_{t \leq T'_i} x_{k_t, t} \geq \alpha_{i-1} \mid \mu_1^K(b^{k'}) \right] \\
 & \geq e^{-\frac{2\alpha_i T'_i}{iK}} \sum_S \left( \prod_{j=1}^{i-1} \alpha_j^{|\{t: x_t = \alpha_j\}|} \prod_{k=1}^K \gamma_k(\bar{b})^{|\{t: k_t = k, x_t = 1\}|} \right) \\
 & = e^{-\frac{2\alpha_i T'_i}{iK}} P \left[ \min_{t \leq T'_i} x_{k_t, t} \geq \alpha_{i-1} \mid \eta_1^K(\bar{b}) \right].
 \end{aligned}$$

the first inequality uses Lemma 17 and the last equality holds for the same reason that Equation 6 holds.  $\square$

In Lemma 15, we turn the bound in Lemma 14 on the probability of obtaining  $\alpha_i$  into a bound on the performance of  $\pi$ . Note that Lemma 15 holds uniformly over the values of  $b_j$  for  $j \neq i$ .

**Lemma 15.** *We have*

$$\frac{1}{K} \sum_{k'=1}^K \mathbb{E} \left[ \min_{t \leq T'_i} x_{k_t, t} \mid \mu_1^K(b^{k'}) \right] \geq 2\alpha_i.$$

*Proof.* We have

$$\begin{aligned}
 & \frac{1}{K} \sum_{k'=1}^K \mathbb{E} \left[ \min_{t \leq T'_i} x_{k_t, t} \mid \mu_1^K(b^{k'}) \right] \\
 & \geq \frac{\alpha_{i-1}}{K} \sum_{k'=1}^K P \left[ \min_{t \leq T'_i} x_{k_t, t} \geq \alpha_{i-1} \mid \mu_1^K(b^{k'}) \right] \quad (10) \\
 & \geq \alpha_{i-1} e^{-\frac{2\alpha_i T'_i}{iK}} P \left[ \min_{t \leq T'_i} x_{k_t, t} \geq \alpha_{i-1} \mid \eta_1^K(\bar{b}) \right]
 \end{aligned}$$

The first inequality is Markov's inequality. The second inequality is Lemma 14. We have

$$\begin{aligned}
 P \left[ \min_{t \leq T'_i} x_{k_t, t} \geq \alpha_{i-1} \mid \eta_1^K(\bar{b}) \right] & \geq \left( 1 - \frac{\alpha_i}{K} - \sum_{j=i+1}^{\infty} \alpha_j \right)^{T'_i} \\
 & \geq \left( 1 - \frac{\alpha_i(1 + \frac{1}{i})}{K} \right)^{T'_i} \\
 & \geq e^{-\alpha_i(1 + \frac{1}{i})^2 T'_i / K} \\
 & \geq \alpha_i^{(1 + \frac{1}{i})^2 c_i}. \quad (11)
 \end{aligned}$$

The first inequality lower bounds the probability of obtaining a value of  $\alpha_i$  or less at every iteration. The second inequality uses Lemma 12(C). The third inequality uses Lemma 18 and Lemma 12(B). The fourth inequality uses the definition  $T'_i = \lfloor c_i K \log(1/\alpha_i)/\alpha_i \rfloor$ .

Combining the Equation 10 and Equation 11 gives

$$\begin{aligned}
 \frac{1}{K} \sum_{k'=1}^K \mathbb{E} \left[ \min_{t \leq T'_i} x_{k_t, t} \mid \mu_1^K(b^{k'}) \right] & \geq \alpha_{i-1} e^{-\frac{2\alpha_i T'_i}{iK}} \alpha_i^{(1 + \frac{1}{i})^2 c_i} \\
 & \geq 2\alpha_i^{\frac{1}{i}} \alpha_i^{\frac{2c_i}{i}} \alpha_i^{(1 + \frac{1}{i})^2 c_i} \\
 & = 2\alpha_i.
 \end{aligned}$$

The second inequality uses Lemma 12(D) and the definition of  $T'_i$ . The third line uses the definition  $c_i = (1 - 1/i)/((1 + 1/i)^2 + 2/i)$ , which was chosen to make the third line hold. This completes the proof of the lemma.  $\square$

Noting that Lemma 15 holds uniformly over the values of  $b_j$  for  $j \neq i$ , a direct consequence of Lemma 15 is Corollary 16.

**Corollary 16.** *We have*

$$P_{b \sim D} \left( \mathbb{E} \left[ \min_{t \leq T'_i} x_{k_t, t} \mid \mu_1^K(b) \right] \geq 2\alpha_i \right) \geq \frac{1}{K},$$

where  $D$  is the distribution over  $\{1, \dots, K\}^\infty$  defined by sampling each component independently and uniformly at random from  $\{1, \dots, K\}$ . The outer probability is over  $b$ , and the inner expectation is over the  $x_{k_t, t}$ .

### 4.3 Proof of Theorem 11

Here we synthesize the above results to prove Theorem 11. Lemma 13 and Corollary 16 together imply that

$$\begin{aligned}
 P_{b \sim D} \left( \mathbb{E} \left[ \min_{t \leq T'_i} x_{k_t, t} \mid \mu_1^K(b) \right] \geq 2\alpha_i > \mathbb{E} \left[ \min_{t \leq T'_i} x_{k_*, t} \mid \mu_1^K(b) \right] \right) \\
 \geq \frac{1}{K},
 \end{aligned}$$

which directly implies that  $P(R_{T'_i}^{\pi, \pi_*} \geq T'_i/T_i) \geq 1/K$ . Recall that for a sequence of events  $A_i$ , we have  $P(\text{infinitely many } A_i \text{ happen}) \geq \limsup P(A_i)$ . This can be seen by applying Fatou's lemma to the relevant indicator functions. It follows that

$$P_{b \sim D} \left( R_{T'_i}^{\pi, \pi_*} \geq \frac{T'_i}{T_i} \text{ for infinitely many } i \right) \geq \frac{1}{K}.$$

Recall the definitions

$$T_i = \lceil \log(1/\alpha_i)/\alpha_i \rceil \quad T'_i = \lfloor c_i K \log(1/\alpha_i)/\alpha_i \rfloor.$$

Since  $c_i \rightarrow 1$ , it follows that  $T'_i/T_i \rightarrow K$ , and so there exists a tuple  $\mu_1^K \in M_K$  such that  $\limsup_T R_T^{\pi, \pi_*} \geq K$ , proving the claim.

## 5 Related Work

Our setting is closely related to the multi-armed bandit problem, which has been studied extensively. See Bubeck and Cesa-Bianchi (2012) for a survey. Regret is the most common measure of performance, though some authors study “simple regret” (Bubeck et al., 2011), where the goal is to identify the arm with the lowest mean. However, these settings provide little guidance on designing a policy to minimize the single smallest cost. The extreme bandit problem, where we care not about the average cost but about the single minimal cost, has been significantly less studied.

The extreme bandit problem (also called the max  $K$ -armed bandit problem) is introduced in Cicirello and Smith (2005) and further developed in Streeter and Smith (2006a,b). The problem is additionally studied in Carpentier and Valko (2014), where the authors give an explicit algorithm and prove that it exhibits asymptotically no regret in the sense of Equation 1. However, all results in previous work have relied heavily on strong parametric or semiparametric assumptions on the distributions  $\mu_1^K$  under consideration. Motivated by extreme value theory, Cicirello and Smith (2005) assume that the distributions belong to the Gumbel family and Carpentier and Valko (2014) consider distributions in the Fréchet family (or distributions that are well approximated by the Fréchet family). When the individual samples arise as the maxima of a large number of independent, identically-distributed random variables, then these assumptions may be realistic. These assumptions dramatically simplify the problem. As in the multi-armed bandit setting, where every sample from a distribution provides information about the mean of the distribution, in the parametric setting, every sample provides information about the parameters of the distribution. Once we have accurately estimated each distribution, we can make sensible choices about which distribution to choose. Our work shows that some form of assumptions are necessary to improve on the guarantees of the policy that chooses each arm equally often.

We do not expect the parametric assumptions motivated by extreme value theory to make sense in the setting of hyperparameter optimization. However, the question of what realistic assumptions are likely to hold in practice for hyperparameter optimization is an important question.

More recently, David and Shimkin (2015) consider a PAC setting for the extreme bandit problem and prove a lower bound on the sample complexity of algorithms that return an answer within  $\epsilon$  of the optimal attainable value with probability  $1 - \delta$ .

The no free lunch theorems are another form of hardness result in the optimization setting. Wolpert and Macready (1997) show that in a discrete setting, all optimization algorithms that never revisit the same point perform equally well in expectation with respect to the uniform distribution over all possible objectives.

## 6 Discussion

We have shown that a number of subtleties arise in the extreme bandit setting that are not present in the standard bandit setting. These include the fact that there is no well-defined “best” arm and the fact that strategies that play multiple arms can outperform oracle strategies that play a single arm. We have shown that no policy can be guaranteed to perform asymptotically as well as an oracle that plays the single best arm for a given time horizon. This result should not be construed to say that no policy can do better in practice. Indeed, hyperparameter optimization problems in the real world possess many nice structural properties. For instance, many hyperparameters have a sweet spot outside of which the algorithm performs poorly. This suggests that many black-box objectives for hyperparameter optimization may exhibit coordinate-wise quasiconvexity. Crafting plausible assumptions on the objectives and understanding how they translate into conditions on the induced distributions over algorithm performance is an important problem.

## Acknowledgements

We would like to thank Balázs Kégl for valuable discussions. We would like to thank Kevin Jamieson and Ilya Tolstikhin for their feedback on earlier drafts of this paper.

## References

- J. Bergstra and Y. Bengio. Random search for hyperparameter optimization. *The Journal of Machine Learning Research*, 13(1):281–305, 2012.
- J. S. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl. Algorithms for hyper-parameter optimization. In *Advances in Neural Information Processing Systems*, pages 2546–2554, 2011.
- S. Bubeck and N. Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.
- S. Bubeck, R. Munos, and G. Stoltz. Pure exploration in finitely-armed and continuous-armed bandits.



bits. *Theoretical Computer Science*, 412(19):1832–1852, 2011.

A. Carpentier and M. Valko. Extreme bandits. In *Advances in Neural Information Processing Systems*, pages 1089–1097, 2014.

V. A. Cicirello and S. F. Smith. The max K-armed bandit: A new model of exploration applied to search heuristic selection. In *Proceedings of the National Conference on Artificial Intelligence*, 2005.

Y. David and N. Shimkin. The max  $k$ -armed bandit: A PAC lower bound and tighter algorithms. *arXiv preprint arXiv:1508.05608*, 2015.

J. Duchi, M. Jordan, M. Wainwright, and A. Wibisono. Optimal rates for zero-order convex optimization: The power of two function evaluations. *IEEE Transactions on Information Theory*, 61(5):2788–2806, 2015.

N. Hansen. The CMA evolution strategy: a comparing review. In *Towards a New Evolutionary Computation. Advances on Estimation of Distribution Algorithms*, pages 75–102. Springer, 2006.

J. C. Lagarias, J. A. Reeds, M. H. Wright, and P. E. Wright. Convergence properties of the Nelder–Mead simplex method in low dimensions. *SIAM Journal on Optimization*, 9(1):112–147, 1998.

M. Powell. The NEWUOA software for unconstrained optimization without derivatives. In *Large-Scale Nonlinear Optimization*, volume 83, pages 255–297. Springer, 2006.

J. Snoek, H. Larochelle, and R. P. Adams. Practical Bayesian optimization of machine learning algorithms. In *Advances in Neural Information Processing Systems*, pages 2951–2959, 2012.

M. J. Streeter and S. F. Smith. An asymptotically optimal algorithm for the max K-armed bandit problem. In *Proceedings of the National Conference on Artificial Intelligence*, 2006a.

M. J. Streeter and S. F. Smith. A simple distribution-free approach to the max K-armed bandit problem. In *Principles and Practice of Constraint Programming-CP 2006*, pages 560–574, 2006b.

Z. Wang, M. Zoghi, F. Hutter, D. Matheson, and N. de Freitas. Bayesian optimization in high dimensions via random embeddings. In *International Joint Conferences on Artificial Intelligence*, 2013.

D. H. Wolpert and W. G. Macready. No free lunch theorems for optimization. *IEEE Transactions on Evolutionary Computation*, 1(1):67–82, 1997.

## A The Best Arm Depends on the Time Horizon

In Example 9, we considered an infinite collection of arms  $\mu_s$  indexed by  $0 < s < 1$ . Samples  $x_{s,t}$  from  $\mu_s$  satisfy  $P(x_{s,t} = s) = s$  and  $P(x_{s,t} = 1) = 1 - s$ . We claimed that for a time horizon of  $T$ , the optimal  $s$  is  $\Theta((\log T)/T)$ .

We have

$$\mathbb{E} \left[ \min_{t \leq T} x_{s,t} \right] = s(1 - (1-s)^T) + 1(1-s)^T = s + (1-s)^{T+1}.$$

Let  $s_*$  be the index of the optimal distribution, so  $\min_s \mathbb{E}[\min_{t \leq T} x_{s,t}] = s_* + (1 - s_*)^{T+1}$ . For large  $T$ , we can consider the range  $0 < s \leq \frac{1}{2}$ . We have

$$s + e^{-2s(T+1)} \leq s + (1-s)^{T+1} \leq s + e^{-s(T+1)}.$$

It follows that

$$\begin{aligned} & s_* + e^{-2s_*(T+1)} \\ & \leq \min_s \mathbb{E} \left[ \min_{t \leq T} x_{s,t} \right] \\ & \leq \min_s s + e^{-s(T+1)} \\ & \leq \frac{\log T}{T+1} + \frac{1}{T} \\ & \leq \frac{2 \log T}{T}. \end{aligned}$$

Therefore,  $s_* \leq (2 \log T)/T$  and  $e^{-2s_*(T+1)} \leq (2 \log T)/T$ . The latter implies that

$$s_* \geq \frac{-\log 2 - \log \log T + \log T}{2(T+1)}$$

These results imply that  $s_*$  is  $\Theta((\log T)/T)$ .

## B Proof of Lemma 17

Here we state and prove Lemma 17, which is used in the proof of Lemma 14. The goal of Lemma 17 is to show that the probability of a particular sequence of values under the tuple  $\mu_1^K(b^{k'})$ , when averaged over the possible values of  $k'$ , is at least as great (up to a constant  $c$ ) as the probability of the same sequence of values under the averaged tuple  $\eta_1^K$ . Since all values other than the values 1 and  $\alpha_i$  have equal probability under all tuples (for  $j \neq i$ , the value  $\alpha_j$  has probability  $\alpha_j$  under the  $b_j$ th element of each tuple), this lemma focuses on the probabilities of the values that equal 1. Recall that  $\gamma_k(b^{k'})$  is the probability of obtaining a value of 1 from  $\mu_k(b^{k'})$  and  $\gamma_k(\bar{b})$  is the probability of obtaining a value of 1 from  $\eta_k$ .

**Lemma 17.** For integers  $n_1, \dots, n_K \geq 0$  such that  $n_k \leq T$ , we have

$$\frac{1}{K} \sum_{k'=1}^K \prod_{k=1}^K \gamma_k(b^{k'})^{n_k} \geq c \prod_{k=1}^K \gamma_k(\bar{b})^{n_k},$$

where  $c = e^{-\frac{2\alpha_i T}{iK}}$ .

*Proof.* This result nearly follows from Jensen's inequality. Indeed, if the function

$$f(c_1, \dots, c_K) = \prod_{k=1}^K \left( 1 - c_k \alpha_i - \sum_{\substack{j=1 \\ j \neq i}}^{\infty} \mathbb{1}[j = k] \alpha_j \right)^{n_k}$$

were convex, then the result would follow from a single application of Jensen's inequality. That is, the result with  $c = 1$  is precisely the statement

$$\frac{f(1, 0, \dots, 0) + \dots + f(0, \dots, 0, 1)}{K} \geq f\left(\frac{1}{K}, \dots, \frac{1}{K}\right).$$

Unfortunately, despite the fact that  $f$  is the product of convex functions (over the relevant domains),  $f$  itself is not convex. To circumvent this difficulty, we will approximate each term with the exponential of an affine function, so that the product of approximations remains convex (because the affine functions simply add). As our approximation is imperfect, we pick up a penalty in the form of the constant  $c$ . Let

$$\omega_k = 1 - \sum_{\substack{j=1 \\ j \neq i}}^{\infty} \mathbb{1}[j = k] \alpha_j \quad \beta_{i,k} = \frac{\alpha_i}{\omega_k},$$

First write

$$\begin{aligned} & \frac{1}{K} \sum_{k'=1}^K \prod_{k=1}^K \gamma_k(b^{k'})^{n_k} \\ &= \frac{1}{K} \sum_{k'=1}^K \prod_{k=1}^K (\omega_k - \mathbb{1}[k' = k] \alpha_i)^{n_k} \\ &= \frac{1}{K} \left( \prod_{k=1}^K \omega_k^{n_k} \right) \sum_{k'=1}^K (1 - \beta_{i,k'})^{n_{k'}}. \end{aligned} \quad (12)$$

Note that by Lemma 12(A), we have  $\omega_{k'} \geq \frac{1}{2}$  and so  $\beta_{i,k'} \leq 2\alpha_i$ . It follows from Lemma 18 and

Lemma 12(B) that we can write

$$\begin{aligned} & \frac{1}{K} \sum_{k'=1}^K (1 - \beta_{i,k'})^{n_{k'}} \\ & \geq \frac{1}{K} \sum_{k'=1}^K e^{-(1+1/i)\beta_{i,k'} n_{k'}} \\ & \geq e^{-(1+1/i)\frac{1}{K} \sum_{k'=1}^K \beta_{i,k'} n_{k'}} \\ & \geq e^{-\frac{2\alpha_i T}{iK}} e^{-\frac{1}{K} \sum_{k'=1}^K \beta_{i,k'} n_{k'}} \\ & \geq e^{-\frac{2\alpha_i T}{iK}} \prod_{k'=1}^K \left( 1 - \frac{\beta_{i,k'}}{K} \right)^{n_{k'}}. \end{aligned} \quad (13)$$

The second inequality is Jensen's inequality. The third inequality breaks the  $1 + 1/i$  term into two terms and uses the bounds  $\beta_{i,k'} \leq 2\alpha_i$  and  $n_{k'} \leq T$ . The fourth inequality uses the fact that  $e^{-x} \geq 1 - x$ . Combining Equation 12 and Equation 13 gives

$$\begin{aligned} & \frac{1}{K} \sum_{k'=1}^K \prod_{k=1}^K \gamma_k(b^{k'})^{n_k} \\ & \geq e^{-\frac{2\alpha_i T}{iK}} \left( \prod_{k=1}^K \omega_k^{n_k} \right) \prod_{k'=1}^K \left( 1 - \frac{\beta_{i,k'}}{K} \right)^{n_{k'}} \\ & = e^{-\frac{2\alpha_i T}{iK}} \prod_{k=1}^K \left( \omega_k - \frac{\alpha_i}{K} \right)^{n_k} \\ & = e^{-\frac{2\alpha_i T}{iK}} \prod_{k=1}^K \gamma_k(\bar{b})^{n_k}, \end{aligned}$$

which finishes the proof.  $\square$

## C Upper Bound on Exponential

Throughout this paper, we make use of the inequality  $e^{-x} \geq 1 - x$ . However, on a couple of occasions, we need to lower bound  $1 - x$  by an exponential of the form  $e^{-rx}$  for some constant  $r$ . The bound that we use is given in Lemma 18.

**Lemma 18.** For  $i \geq 1$  and  $y \in [0, \frac{1}{2(1+i)}]$ , we have  $e^{-y(1+\frac{1}{i})} \leq 1 - y$ .

*Proof.* More generally, the convexity of  $e^{-x}$  implies that for  $0 \leq x \leq c$ , we have

$$e^{-x} \leq 1 - \frac{1 - e^{-c}}{c} x.$$

The right hand side is the formula for the line interpolating between the points  $(0, 1)$  and  $(c, e^{-c})$  on the graph of  $e^{-x}$ . Choosing  $c = \log(1 + \frac{1}{i})$ , and noting that  $0 \leq x \leq \frac{1}{1+i}$  implies that  $0 \leq x \leq c$  because

of the standard inequality  $1 - \frac{1}{x} \leq \log x$ , we see that  $0 \leq x \leq \frac{1}{1+i}$  implies that

$$e^{-x} \leq 1 - \frac{1 - \frac{i}{1+i}}{\log(1 + \frac{1}{i})} x \leq 1 - \frac{\frac{1}{1+i}}{\frac{1}{i}} x = 1 - \frac{i}{1+i} x.$$

Setting  $y = \frac{i}{1+i} x$  and using the fact that  $\frac{1}{2(1+i)} \leq \frac{i}{(1+i)^2}$  gives the result.  $\square$